

# Optimised Auto-scaling for Cloud-based Web Service



Jing Jiang

FACULTY OF ENGINEERING AND INFORMATION TECHNOLOGY

UNIVERSITY OF TECHNOLOGY SYDNEY

A thesis submitted for the degree of

*Doctor of Philosophy*

February 2015

## **CERTIFICATE OF AUTHORSHIP/ORIGINALITY**

I certify that the work in this thesis has not previously been submitted for a degree nor has it been submitted as part of requirements for a degree except as fully acknowledged within the text.

I also certify that the thesis has been written by me. Any help that I have received in my research work and the preparation of the thesis itself has been acknowledged. In addition, I certify that all information sources and literature used are indicated in the thesis.

**Signature of Candidate**

---

## Acknowledgements

This thesis is the fruit of my hard work during the whole of my PhD study since being in Australia. It has not been a straight path but it has been an amazing journey in which I have experienced happiness, excitement, depression, struggle and friendship. I feel quite lucky that so many people have offered me a helping hand in this journey. Now is the right time to thank all of them!

First, I would like to express my sincere appreciation and great gratitude to my principal supervisor, Professor Jie Lu, and my co-supervisor, Professor Guangquan Zhang. They are not only mentors in my research life but also the guides in my everyday life, especially in my first year in Australia. In my research work, they offered me ample trust and free space to choose the topic and never doubted my progress, but always made valuable suggestions if I encountered problems. They also guided me on how to express my research results efficiently including how to organise a paper, the right steps to take in writing, and even how to express ideas clearly in sentences. They both always work so

hard and have a strict attitude to research work, which stimulates me to further efforts and from which I will benefit throughout my life. Besides the research work, they were enormously helpful and caring and never said “No” to my requests - for example, they drove to pick up me at the airport on the first day I arrived in Sydney; made suggestions for renting a house; taught me to cook; invited me to their home for celebrations, and even guided me in managing work and life. There are no words to express all my thanks to them; the only way to repay them was to work hard.

I kindly thank all my colleagues in the Decision Systems and e-Service Intelligent (DeSI) lab for their careful participation in my presentation and valuable comments for my research. I also thank many members of the Centre for Quantum Computation and Intelligent Systems (QCIS) and the Faculty of Engineering and Information Technology, University of Technology, Sydney (UTS) for their various assistance and advice which have been of great benefit to this study. I appreciate the financial support from the International Postgraduate Research Scholarship (IPRS) and UTS President’s Scholarship (UTSP).

I appreciate the conference and research travel support by Prof Chengqi Zhang, director of QCIS, the UTS Vice Chancellor’s conference fund,

and the Faculty of Engineering and Information Technology travel fund.

I sincerely thank Ms. Sue Felix and Ms. Barbara Munday for polishing my English language and doing proofreading for my publications and this thesis; their comments have greatly helped me to improve my academic writing. I would also sincerely thank Ms. Teraesa Ashworth who spent much time introducing me to local culture and helping to improve my spoken English.

I would like to express my earnest appreciation and gratitude to my parents for all their selfless love and generous understanding, also to my brothers and sisters for their warm concern and for looking after our parents when I was not around for all those years.

Last but certainly not least, I feel extremely fortunate that I could share the joy and the pain with my dear husband, Guodong Long. He has been by my side during hard times and has comforted me and encouraged me to continue. I am grateful for his love and endless optimism at times when the barriers seemed impossible to break down.

# **Abstract**

Elasticity and cost-effectiveness are two key features for ensuring that cloud-based web services appeal to more businesses. However, true elasticity and cost-effectiveness in the pay-per-use cloud business model has not yet been fully achieved. The explosion of cloud-based web services brings new challenges to enable the automatic scaling up and down of service provision when the workload is time-varying.

This research studies the problems associated with these challenges. It proposes a novel scheme to achieve optimised auto-scaling for cloud-based web services from three levels of cloud structure: Software as a Service (SaaS), Platform as a Service (PaaS), and Infrastructure as a Service (IaaS). At the various levels, auto-scaling for cloud-based web services has different problems and requires different solutions.

At the SaaS level, this study investigates how to design and develop scalable web services, especially for time-consuming applications. To achieve the greatest efficiency, the optimisation of service provision problem is studied by providing the minimum functionality and fastest

scalability performance concerning the speed-up curve and QoS (Quality of Service) of the SLA (Service-Level Agreement). At the PaaS level, this work studies how to support dynamic re-configuration when workloads change and the effective deployment of various kinds of web services to the cloud. To achieve optimised auto-scaling of this deployment, a platform is designed to deploy all web services automatically with the minimal number of cloud resources by satisfying the QoS of SLAs. At the IaaS level for two infrastructure resources of virtual machine (VM) and virtual network (VN), this research focuses on studying two types of cloud-based web service: computation-intensive and bandwidth-intensive. To address the optimised auto-scaling problem for computation-intensive cloud-based web service, data-driven VM auto-scaling approaches are proposed to handle the workload in both stable and dynamic environments. To address the optimised auto-scaling problem for bandwidth-intensive cloud-based web service, this study proposes a novel approach to predict the volume of requests and dynamically adjust the software defined network (SDN)-based network configuration in the cloud to auto-scale the service with minimal cost.

This research proposes comprehensive and profound perspectives to

solve the auto-scaling optimisation problems for cloud-based web services. The proposed approaches not only enable cloud-based web services to minimise resource consumption while auto-scaling service provision to achieve satisfying performance, but also save energy consumption for the global realisation of green computing. The performance of the proposed approaches has been evaluated on a public platform (e.g. Amazon EC2) with the real dataset workload of web services. The experiment results demonstrate that the proposed approaches are practicable and achieve superior performance to other benchmark methods.



# Table of Contents

<b>Table of Contents</b>	<b>viii</b>
<b>List of Figures</b>	<b>xii</b>
<b>List of Tables</b>	<b>xiv</b>
<b>1 Introduction</b>	<b>1</b>
1.1 Research Questions . . . . .	4
1.2 Research Objectives . . . . .	6
1.3 Research Significance . . . . .	9
1.4 Research Methodology . . . . .	10
1.5 Thesis Structure . . . . .	13
1.6 Publications Related to the Thesis . . . . .	16
<b>2 Literature Review</b>	<b>18</b>
2.1 Auto-scaling Cloud-based Web Service (ACSW) . . . . .	18

## TABLE OF CONTENTS

---

2.1.1	Cloud Computing . . . . .	18
2.1.2	Cloud Resource Scaling . . . . .	21
2.1.3	Auto-Scaling for Cloud Resource . . . . .	23
2.2	Optimised Auto-scaling on Cloud Computing . . . . .	30
2.2.1	SaaS Level Optimised Auto-scaling for Cloud Computing .	31
2.2.2	PaaS Level Optimised Auto-scaling for Cloud Computing .	33
2.2.3	IaaS Level Optimised Auto-scaling for Cloud Computing .	36
<b>3</b>	<b>ACWS Optimisation on SaaS</b>	<b>43</b>
3.1	Methodology . . . . .	45
3.1.1	Scaling-up ICF Algorithm on MapReduce . . . . .	50
3.1.2	Computing Average Rating for Items . . . . .	51
3.1.3	Computing Similarity . . . . .	53
3.1.4	Computing Prediction Matrix . . . . .	55
3.1.5	Optimisation Algorithm . . . . .	58
3.2	Experimental Evaluation . . . . .	61
3.2.1	Experiment Setup . . . . .	61
3.2.2	Performance Evaluation . . . . .	62
3.2.3	Summary . . . . .	65
<b>4</b>	<b>ACWS Optimisation on PaaS</b>	<b>66</b>
4.1	Methodology . . . . .	68

## TABLE OF CONTENTS

---

4.1.1	Architecture of Self-adaptive Configuration Optimisation	
	System . . . . .	69
4.1.2	SLA Metrics Modelling using Utility Function . . . . .	72
4.1.3	Stochastic SLA Metrics Computation . . . . .	74
4.1.4	Optimisation Algorithm . . . . .	77
4.2	Experimental Evaluation . . . . .	79
4.2.1	Experiment Setup . . . . .	79
4.2.2	Performance Evaluation . . . . .	81
4.3	Summary . . . . .	85
<b>5</b>	<b>ACWS Optimisation on IaaS for VMs</b>	<b>87</b>
5.1	Methodology . . . . .	91
5.1.1	Overview of the Method . . . . .	91
5.1.2	Prediction Model . . . . .	93
5.1.3	Optimisation Model . . . . .	99
5.2	Experimental Evaluation . . . . .	102
5.2.1	Experiment Setup and Datasets . . . . .	103
5.2.2	Features Selection Evaluation . . . . .	104
5.2.3	Evaluation Methods . . . . .	105
5.2.4	Prediction Model Evaluation . . . . .	107
5.2.5	Allocation Evaluation . . . . .	108

## TABLE OF CONTENTS

---

5.2.6	Performance Evaluation for a Web Application . . . . .	109
5.3	Summary . . . . .	110
<b>6</b>	<b>ACWS Optimisation on IaaS for Bandwidth</b>	<b>113</b>
6.1	Methodology . . . . .	117
6.1.1	Dynamic Prediction Model . . . . .	118
6.1.2	VN Consumption Model . . . . .	120
6.1.3	Profit Maximisation of VN Consumption . . . . .	125
6.2	Experimental Evaluation . . . . .	129
6.2.1	Experiment Setup and Datasets . . . . .	130
6.2.2	Evaluation of Neural Network-based Prediction Model . . . . .	131
6.2.3	Evaluation of Network Consumption Auto-scaling . . . . .	133
6.2.4	Performance Study on Various Parameters . . . . .	136
6.3	Summary . . . . .	137
<b>7</b>	<b>Conclusions and Future Research</b>	<b>138</b>
7.1	Conclusions . . . . .	138
7.2	Main Contributions . . . . .	139
7.3	Limitations and Further Studies . . . . .	140
	<b>Abbreviations</b>	<b>144</b>
	<b>References</b>	<b>148</b>

# List of Figures

1.1	Thesis Structure . . . . .	14
2.1	Summary of Cloud Scaling Techniques . . . . .	32
3.1	MapReduce Computation Framework . . . . .	47
3.2	Workflow of Item-based Collaborative Filtering Algorithm . . . . .	51
3.3	Speedup and Number of Nodes . . . . .	59
3.4	Speedup of Item-based CF Algorithm . . . . .	63
3.5	Isoefficiency Function for Varying Nodes and Data-size with Fixed Running-time . . . . .	64
4.1	An Example of Deployment Configuration . . . . .	69
4.2	Architecture of Self-adaptive Configuration System in Cloud Com- puting . . . . .	70
4.3	Self-adaptive Reconfiguration Interval . . . . .	70
4.4	Request Scaling for Time-varying Workload . . . . .	80

## LIST OF FIGURES

---

4.5	Cost for Time-varying Workload . . . . .	82
4.6	Utilisation for Time-varying Workload . . . . .	83
4.7	Latency for Time-varying Workload . . . . .	84
4.8	Throughput for Time-varying Workload . . . . .	85
5.1	Overview of Optimised Cloud Resource Auto-scaling . . . . .	92
5.2	Transition Rate for the Web Requests Process on VMs . . . . .	96
5.3	Prediction and Allocation with a Dynamic Cap . . . . .	108
5.4	Allocation with Queuing Theory. . . . .	109
5.5	Allocation Comparison for Different Methods . . . . .	110
6.1	Execution Paradigm of SDN-based Bandwidth Auto-scaling . . . .	118
6.2	Transition Rate for the Web Requests Process on Channels . . . .	123
6.3	Rate Transition Rate for the Web Requests Process on Channels .	126
6.4	The Prediction Curve for DS-3 . . . . .	134
6.5	Performance Impact on Various Parameters with DS-3 . . . . .	136
7.1	Summary of Cloud Scaling Techniques . . . . .	141

# List of Tables

5.1	Average SKL Divergence on Different Period Vectors . . . . .	105
5.2	Top 10 Correlated Lags . . . . .	105
5.3	Performance of Regression Model . . . . .	106
5.4	The Confidence Interval with Different Paddings . . . . .	107
5.5	Performance Comparison for Different Methods . . . . .	111
6.1	Performance of Regression Method . . . . .	133
6.2	Performance of Regression Model . . . . .	135